

CPSC 330

Applied Machine Learning

Lecture 21: Ethics

UBC 2022-23

Instructor: Mathias Lécuyer

Imports

```
In [29]: 1 import os
          2 import sys
          3
          4 import IPython
          5 import matplotlib.pyplot as plt
          6 #import mglearn
          7 import numpy as np
          8 import pandas as pd
          9 from IPython.display import HTML, display
         10 from sklearn.dummy import DummyClassifier
         11 from sklearn.linear_model import LogisticRegression
         12 from sklearn.metrics import accuracy_score, f1_score, precision_score,
         13 from sklearn.model_selection import cross_val_score, cross_validate, tr
         14 from sklearn.pipeline import Pipeline, make_pipeline
         15 from sklearn.preprocessing import StandardScaler
         16
         17 %matplotlib inline
         18 pd.set_option("display.max_colwidth", 200)
         19
         20 from IPython.display import Image
```

Lecture plan

- Guest lecture by Dr. Giulia Toti!
- ML fairness activity

ML fairness activity

AI/ML systems can give the illusion of objectivity as they are derived from seemingly unbiased data & algorithm. However, human are inherently biased and AI/ML systems, if not carefully evaluated, can even further amplify the existing inequities and systemic bias in our society.

How do we make sure our AI/ML systems are *fair*? Which metrics can we use to quatify 'fairness' in AI/ML systems?

Let's examine this on [the adult census data set \(https://www.kaggle.com/uciml/adult-census-income\)](https://www.kaggle.com/uciml/adult-census-income).

```
In [30]: 1 census_df = pd.read_csv("../data/adult.csv")
          2 census_df.shape
```

Out[30]: (32561, 15)

```
In [31]: 1 train_df, test_df = train_test_split(census_df, test_size=0.4, random_s
```

```
In [32]: 1 train_df
```

```
Out[32]:
```

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship
25823	36	Private	245521	7th-8th	4	Married-civ-spouse	Farming-fishing	Husband
10274	26	Private	134287	Assoc-voc	11	Never-married	Sales	Own-child
27652	25	Local-gov	109526	HS-grad	9	Married-civ-spouse	Craft-repair	Husband
13941	23	Private	131275	HS-grad	9	Never-married	Craft-repair	Own-child
31384	27	Private	193122	HS-grad	9	Married-civ-spouse	Machine-op-inspct	Husband
...
29802	25	Private	410240	HS-grad	9	Never-married	Craft-repair	Own-child
5390	51	Private	146767	Assoc-voc	11	Married-civ-spouse	Prof-specialty	Husband
860	55	Federal-gov	238192	HS-grad	9	Married-civ-spouse	Tech-support	Husband
15795	41	Private	154076	Some-college	10	Married-civ-spouse	Adm-clerical	Husband
23654	22	Private	162667	HS-grad	9	Never-married	Handlers-cleaners	Own-child

19536 rows × 15 columns

```
In [33]: 1 train_df_nan = train_df.replace("?", np.nan)
          2 test_df_nan = test_df.replace("?", np.nan)
          3 train_df_nan.shape
```

```
Out[33]: (19536, 15)
```

```
In [34]: 1 # Let's identify numeric and categorical features
          2
          3 numeric_features = [
          4     "age",
          5     "fnlwgt",
          6     "capital.gain",
          7     "capital.loss",
          8     "hours.per.week",
          9 ]
         10
         11 categorical_features = [
         12     "workclass",
         13     "marital.status",
         14     "occupation",
         15     "relationship",
         16     # "race",
         17     "native.country",
         18 ]
         19
         20 ordinal_features = ["education"]
         21 binary_features = [
         22     "sex"
         23 ] # Not binary in general but in this particular dataset it seems to h
         24 drop_features = ["education.num"]
         25 target = "income"
```

```
In [35]: 1 train_df["education"].unique()
```

```
Out[35]: array(['7th-8th', 'Assoc-voc', 'HS-grad', 'Bachelors', 'Some-college',
                '10th', '11th', 'Prof-school', '12th', '5th-6th', 'Masters',
                'Assoc-acdm', '9th', 'Doctorate', '1st-4th', 'Preschool'],
              dtype=object)
```

```
In [36]: 1 education_levels = [  
2     "Preschool",  
3     "1st-4th",  
4     "5th-6th",  
5     "7th-8th",  
6     "9th",  
7     "10th",  
8     "11th",  
9     "12th",  
10    "HS-grad",  
11    "Prof-school",  
12    "Assoc-voc",  
13    "Assoc-acdm",  
14    "Some-college",  
15    "Bachelors",  
16    "Masters",  
17    "Doctorate",  
18    ]
```

```
In [37]: 1 assert set(education_levels) == set(train_df["education"].unique())
```

```
In [38]: 1 X_train = train_df_nan.drop(columns=[target])  
2 y_train = train_df_nan[target]  
3  
4 X_test = test_df_nan.drop(columns=[target])  
5 y_test = test_df_nan[target]
```

```
In [39]: 1 from sklearn.compose import ColumnTransformer, make_column_transformer  
2 from sklearn.impute import SimpleImputer  
3 from sklearn.preprocessing import OneHotEncoder, OrdinalEncoder, StandardScaler  
4  
5 numeric_transformer = make_pipeline(StandardScaler())  
6  
7 ordinal_transformer = OrdinalEncoder(categories=[education_levels], dtype=int)  
8  
9 categorical_transformer = make_pipeline(  
10     SimpleImputer(strategy="constant", fill_value="missing"),  
11     OneHotEncoder(handle_unknown="ignore", sparse=False),  
12 )  
13  
14 binary_transformer = make_pipeline(  
15     SimpleImputer(strategy="constant", fill_value="missing"),  
16     OneHotEncoder(drop="if_binary", dtype=int),  
17 )  
18  
19 preprocessor = make_column_transformer(  
20     (numeric_transformer, numeric_features),  
21     (ordinal_transformer, ordinal_features),  
22     (binary_transformer, binary_features),  
23     (categorical_transformer, categorical_features),  
24     ("drop", drop_features),  
25 )
```

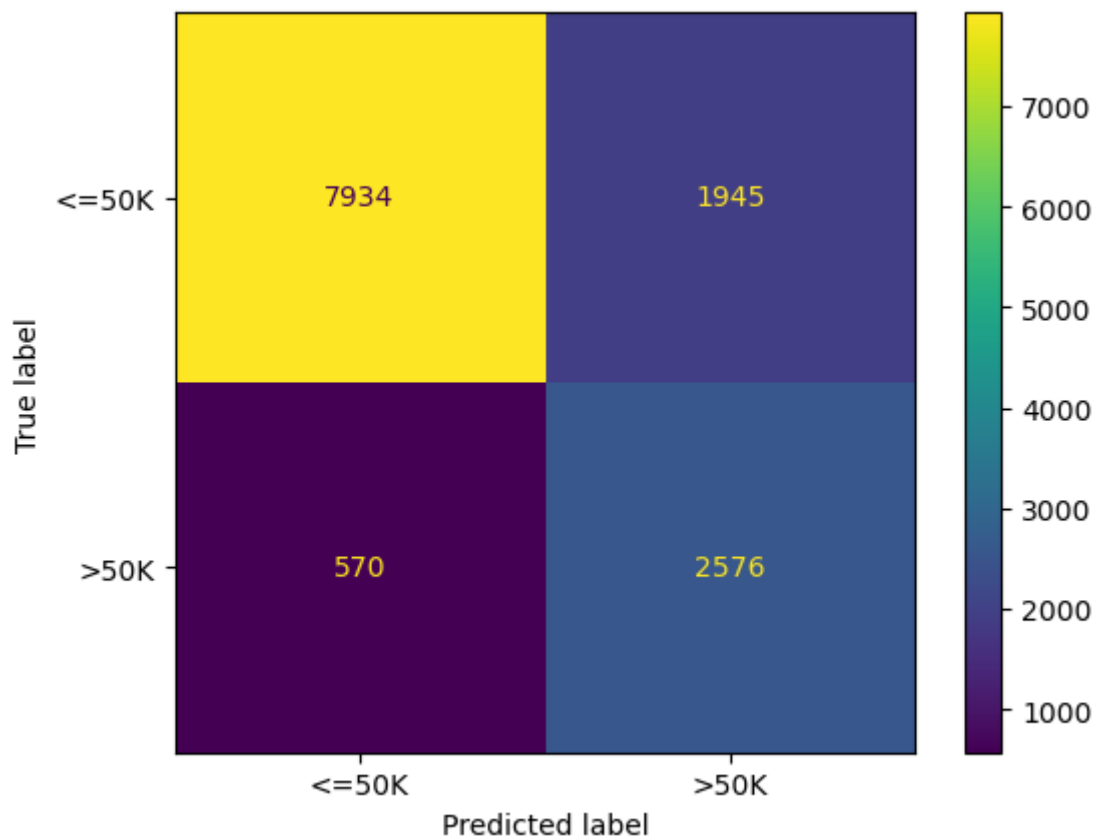
```
In [40]: 1 y_train.value_counts()
```

```
Out[40]: <=50K      14841
         >50K       4695
         Name: income, dtype: int64
```

```
In [41]: 1 pipe_lr = make_pipeline(
         2     processor, LogisticRegression(class_weight="balanced", max_iter=
         3 )
```

```
In [42]: 1 pipe_lr.fit(X_train, y_train);
```

```
In [43]: 1 from sklearn.metrics import ConfusionMatrixDisplay # Recommended metho
         2
         3 ConfusionMatrixDisplay.from_estimator(pipe_lr, X_test, y_test);
```



Let's examine confusion matrix separately for the two genders we have in the data.

```
In [44]: 1 X_train_enc = processor.fit_transform(X_train)
         2 processor.named_transformers_["pipeline-2"]["onehotencoder"].get_fea
```

```
Out[44]: array(['x0_Male'], dtype=object)
```

In [45]: 1 X_test.head()

Out[45]:

	age	workclass	fnlwgt	education	education.num	marital.status	occupation	relationship	
14160	29	Private	280618	Some-college	10	Married-civ-spouse	Handlers-cleaners	Husband	✓
27048	19	Private	439779	Some-college	10	Never-married	Sales	Own-child	✓
28868	28	Private	204734	Some-college	10	Married-civ-spouse	Tech-support	Wife	✓
5667	35	Private	107991	11th	7	Never-married	Sales	Not-in-family	✓
7827	20	Private	54152	Some-college	10	Never-married	Adm-clerical	Own-child	✓

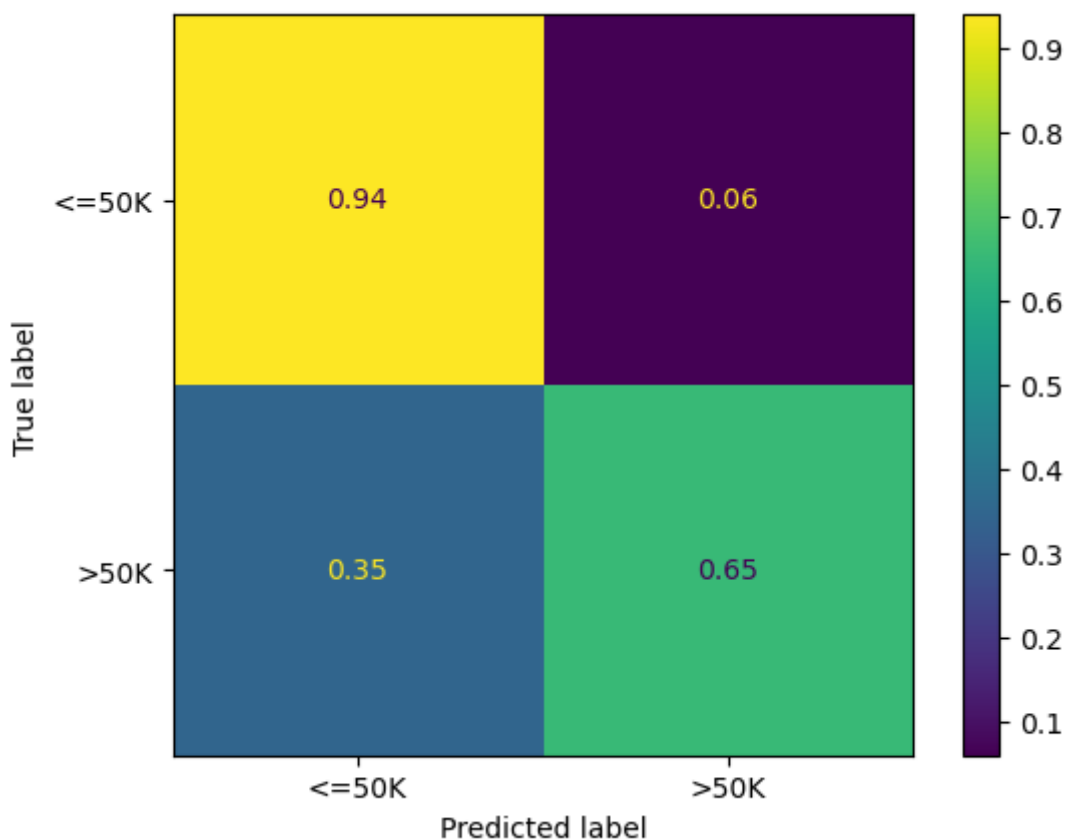
In [46]:

```

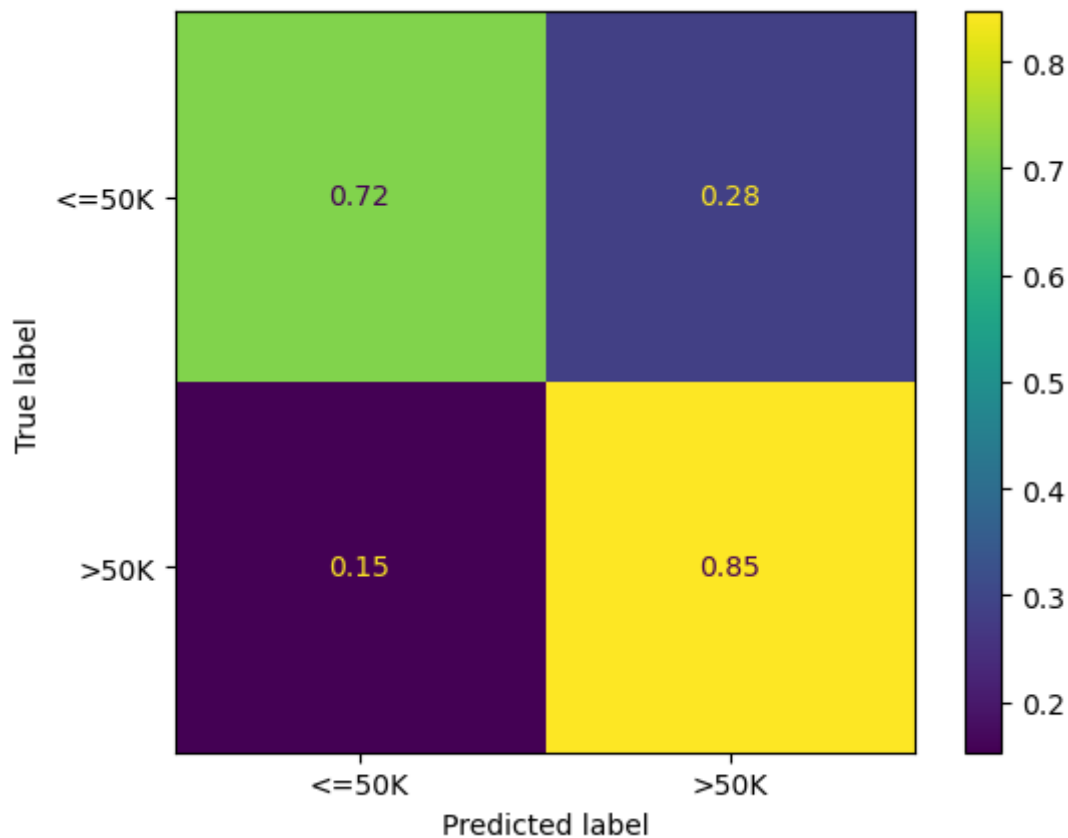
1 X_female = X_test.query("sex=='Female'")
2 X_male = X_test.query("sex=='Male'")
3
4 y_female = y_test[X_female.index]
5 y_male = y_test[X_male.index]
6 female_preds = pipe_lr.predict(X_female)
7 male_preds = pipe_lr.predict(X_male)

```

In [47]: 1 ConfusionMatrixDisplay.from_estimator(pipe_lr, X_female, y_female, norm



```
In [48]: 1 ConfusionMatrixDisplay.from_estimator(pipe_lr, X_male, y_male, normaliz
```



What's the accuracy of this model?

```
In [49]: 1 from sklearn.metrics import confusion_matrix
2
3 data = {"male": [], "female": []}
4 f_TN, f_FP, f_FN, f_TP = confusion_matrix(y_female, female_preds).ravel()
5 m_TN, m_FP, m_FN, m_TP = confusion_matrix(y_male, male_preds).ravel()
```

```
In [50]: 1 accuracy_male = accuracy_score(y_male, male_preds)
2 accuracy_female = accuracy_score(y_female, female_preds)
3 data["male"].append(accuracy_male)
4 data["female"].append(accuracy_female)
5
6 print("Accuracy male: {:.3f}".format(accuracy_male))
7 print("Accuracy female: {:.3f}".format(accuracy_female))
```

Accuracy male: 0.756
Accuracy female: 0.909

```
In [ ]: 1
```

```
In [51]: 1 y_female.value_counts(normalize=True)
```

```
Out[51]: <=50K    0.892675
>50K      0.107325
Name: income, dtype: float64
```

```
In [52]: 1 y_male.value_counts(normalize=True)
```

```
Out[52]: <=50K    0.691999
>50K      0.308001
Name: income, dtype: float64
```

There is more class imbalance for female!

Let's assume that a company is using this classifier for loan approval with a simple rule that if the income is $\geq 50K$, approve the loan else reject the loan.

Statistical parity suggests that the proportion of each segment of a protected class (e.g. sex) should receive the positive outcome at equal rates. For example, the number of loans approved for female should be equal to male.

Calculate the precision for male and female. Based on your results, do you think this income classifier is fair?

```
In [53]: 1 precision_male = precision_score(y_male, male_preds, pos_label=">50K")
2 precision_female = precision_score(y_female, female_preds, pos_label=">50K")
3 data["male"].append(precision_male)
4 data["female"].append(precision_female)
5
6 print("Precision male: {:.3f}".format(precision_male))
7 print("Precision female: {:.3f}".format(precision_female))
```

```
Precision male: 0.570
Precision female: 0.567
```

Equal opportunity suggests that each group should get the positive outcome at equal rates, assuming that people in this group qualify for it. For example, if a man and a woman have both a certain level of income, we want them to have the same chance of getting the loan. In other words, the true positive rate (TPR or recall) of both groups should be equal.

```
In [54]: 1 recall_male = recall_score(y_male, male_preds, pos_label=">50K")
2 recall_female = recall_score(y_female, female_preds, pos_label=">50K")
3
4 data["male"].append(recall_male)
5 data["female"].append(recall_female)
6
7 print("Recall male: {:.3f}".format(recall_male))
8 print("Recall female: {:.3f}".format(recall_female))
```

```
Recall male: 0.847
Recall female: 0.654
```


There is usually a tradeoff between rationality (adopting effective means to achieve your desired outcome) and bias. The desired outcome of banks, for example, is to maximize their profit. So in many circumstances, they not only care about approving as many qualified applications as possible (true positive), but also to avoid approving unqualified applications (false positive) because default loan could have detrimental effects for them.

Let's examine false positive rate (FPR) of both groups.

```
In [55]: 1 fpr_male = m_FP / (m_FP + m_TN)
          2 fpr_female = f_FP / (f_FP + f_TN)
          3
          4 data["male"].append(fpr_male)
          5 data["female"].append(fpr_female)
          6
          7 print("FPR male: {:.3f}".format(fpr_male))
          8 print("FPR female: {:.3f}".format(fpr_female))
```

```
FPR male: 0.284
FPR female: 0.060
```

```
In [56]: 1 pd.DataFrame(data, index=["accuracy", "precision", "recall", "FPR"])
```

```
Out[56]:
```

	male	female
accuracy	0.756170	0.909365
precision	0.570103	0.567416
recall	0.847186	0.654428
FPR	0.284340	0.059984

- Discuss these results with your neighbours.
- Does the effect still exist if the sex feature is removed from the model (but you still have it available separately to do the two confusion matrices)?

```
In [ ]: 1
```

```
In [ ]: 1
```